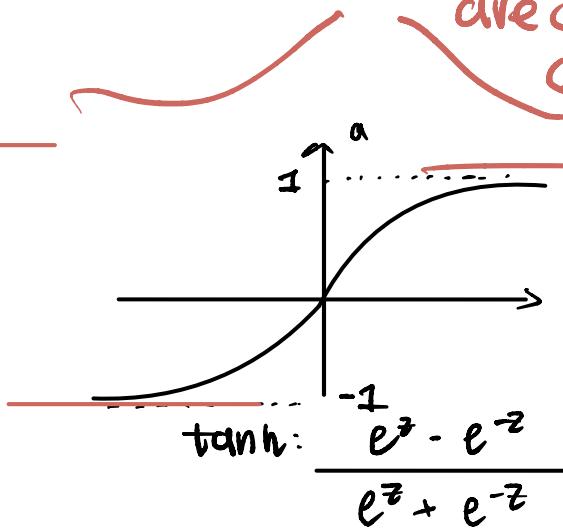
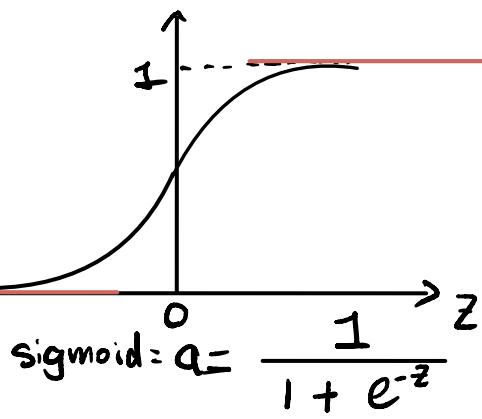
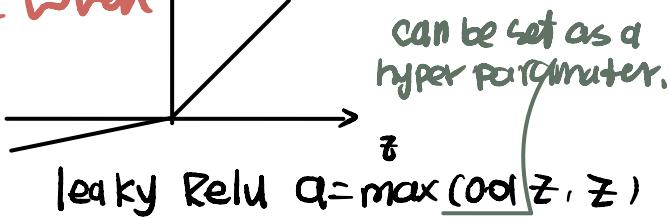
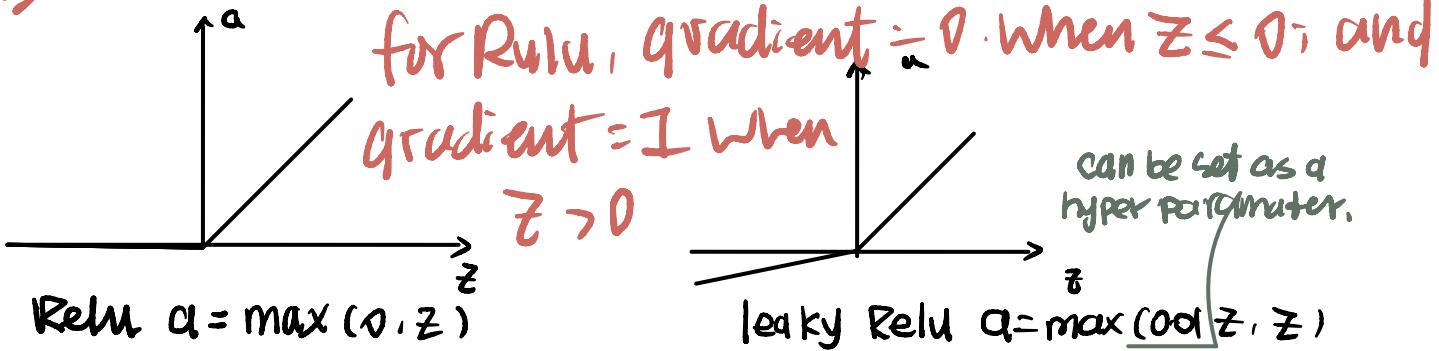


Activation Functions.

Slopes at these points are close to zero which can slow down gradient descent.



⇒ Solution: ReLU and leaky ReLU functions:



why do you need non-linear Activation functions?

counter example: only use linear Activation functions.

Given x :

$$z^{[1]} = W^{[1]}x + b^{[1]} \quad (\text{identity})$$

$a^{[1]} = z^{[1]}$ $g(z) = z$: linear activation function "no use of non-linear activation function here."

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = z^{[2]}$$

$$\rightarrow q^{[1]} = z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[2]} = z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$= (w^{[1]}x + b^{[1]})w^{[2]} + b^{[2]}$$

$$= \frac{(w^{[1]}w^{[2]})x}{w'} + \frac{(b^{[1]}w^{[2]} + b^{[2]})}{b'}$$

$$= w'x + b' \text{ still a linear function}$$

\rightarrow the hidden layers of the neural net
are redundant.

Derivatives of non-linear activation functions.

Sigmoid:

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = \frac{a(1-a)}{\text{activation}}$$

$$a = g(z)$$

Tanh:

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - a^2$$

Relu:

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Although $g(z)$ at $z=0$ is actually not defined.

Leaky Relu:

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Tanh is universally better than Sigmoid, except the case where we are dealing a regression problem. Relu and leaky Relu, on the other hand, are almost always more efficient. Relu is most often set as default activation function.